# Global Convergence of Least Squares EM for Demixing Two Log-Concave Densities

Wei Qian, Yuqian Zhang, Yudong Chen  {wq34,yz2557,yudong.chen}@cornell.edu

**Cornell**Engineering
Operations Research and Information Engineering

## Motivation

The **Expectation Maximization** (EM) algorithm has only a few theoretical guarantees for convergence despite its popularity.

► Recent progress on global convergence has focused on a balanced mixture of 2 Gaussian distributions.

► Can we develop a global convergence theory for a mixture of some broader class of distributions?

## Problem Overview

► **Distribution Class**: log-concave and rotation invariant
$$\mathcal{F} = \left\{ f : f(\boldsymbol{x}) = \frac{1}{C_g} \exp\left( -g(\|\boldsymbol{x}\|_2) \right), \right.$$
$$g \text{ convex and increasing on } [0, \infty),$$
$$\left. \int f(\boldsymbol{x}) \, d\boldsymbol{x} = 1, \int x_i^2 f(\boldsymbol{x}) \, d\boldsymbol{x} = 1, \forall i \in [d] \right\}.$$
Each $f \in \mathcal{F}$ generates a location-scale family consisting of the densities
$$f_{\beta,\sigma}(\boldsymbol{x}) := \frac{1}{\sigma^d} f\left( \frac{\boldsymbol{x} - \beta}{\sigma} \right)$$

► **A Balanced 2-Mixture Generative Model**
$$D(\beta^*, \sigma) := \frac{1}{2} f_{\beta^*,\sigma} + \frac{1}{2} f_{-\beta^*,\sigma}.$$

► **Location Estimation Problem**: Given data $\boldsymbol{X}^1, \ldots, \boldsymbol{X}^n \in \mathbb{R}^d$ sampled i.i.d. from the mixture distribution $D(\beta^*, \sigma)$, $\sigma$ is known, how to estimate $\beta^*$?

► **Classical EM** does not have a closed-form solution for the M-step.

► We analyze the **Least-Squares EM** (LS-EM):

  ► **E-step:** Compute the conditional probabilities given $\beta$,
  $$p^1_{\beta,\sigma}(\boldsymbol{X}) := \frac{f_{\beta,\sigma}(\boldsymbol{X})}{f_{\beta,\sigma}(\boldsymbol{X}) + f_{-\beta,\sigma}(\boldsymbol{X})};$$
  $$p^2_{\beta,\sigma}(\boldsymbol{X}) := \frac{f_{-\beta,\sigma}(\boldsymbol{X})}{f_{\beta,\sigma}(\boldsymbol{X}) + f_{-\beta,\sigma}(\boldsymbol{X})}.$$

  ► <span style="color:red">Least-Squares M-step:</span> weighted least squares regression
  $$M(\beta^*, \beta) = \underset{\boldsymbol{b}}{\operatorname{argmin}} \, \mathbb{E}_{\boldsymbol{X} \sim D(\beta^*,\sigma)} \left[ p^1_{\beta,\sigma}(\boldsymbol{X}) \|\boldsymbol{X} - \boldsymbol{b}\|_2^2 + p^2_{\beta,\sigma}(\boldsymbol{X}) \|\boldsymbol{X} + \boldsymbol{b}\|_2^2 \right]$$
  $$= \mathbb{E}_{\boldsymbol{X} \sim D(\beta^*,\sigma)} \boldsymbol{X} \tanh\left( \frac{1}{2} g\left( \frac{1}{\sigma} \|\boldsymbol{X} + \beta\|_2 \right) - \frac{1}{2} g\left( \frac{1}{\sigma} \|\boldsymbol{X} - \beta\|_2 \right) \right).$$

## Properties of Least Squares EM iterates

► **Two Dimensional Structure**: The LS-EM iterate $M(\beta^*, \beta)$ is in the span of $\beta$ and $\beta^*$.

  ► Invariant 1-dim subspace: in the direction of $\beta^*$ or in the orthogonal direction to $\beta^*$.

► **Angle Decreasing Property**: The angle between the LS-EM iterate $M(\beta^*, \beta)$ and $\text{sign}(\langle \beta, \beta^* \rangle)\beta^*$ is smaller than the angle between $\beta$ and $\text{sign}(\langle \beta, \beta^* \rangle)\beta^*$ when $\beta$ is not orthogonal to $\beta^*$.
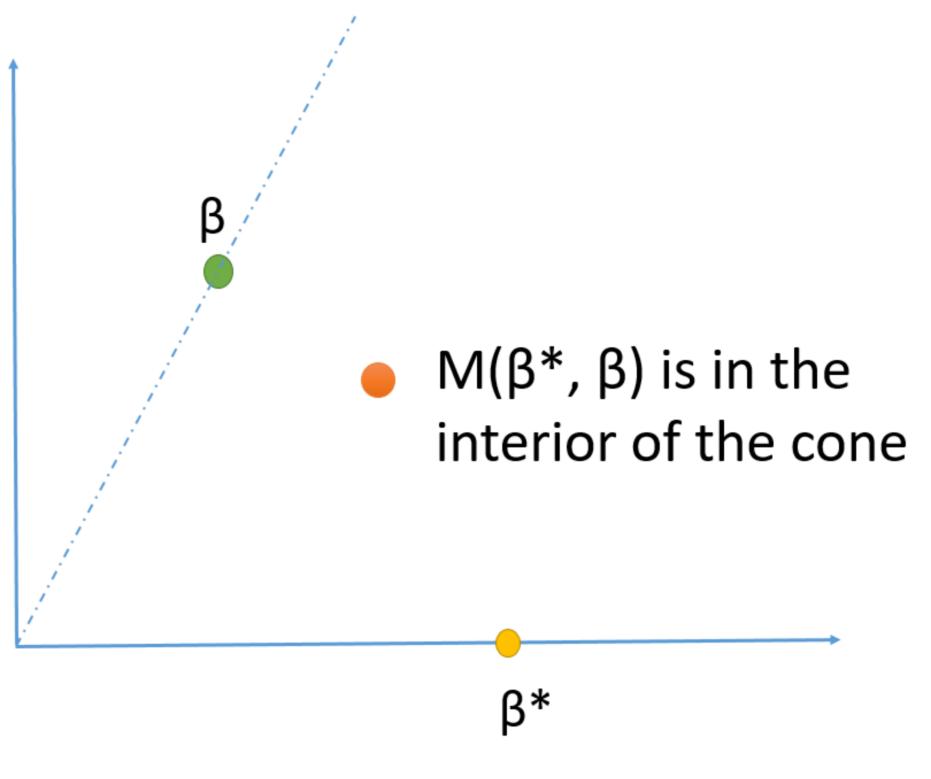


β

● M(β*, β) is in the interior of the cone

β*

Figure: The LS-EM iterate has a smaller angle with $\beta^*$.

► **Asymptotic Convergence**: The Least-Squares EM algorithm converges to $\text{sign}(\langle \beta^0, \beta^* \rangle)\beta^*$ from any randomly initialized point $\beta^0$ that is not orthogonal to $\beta^*$

  ► The angle decreasing property forces the iterates to converge to the correct subspace;

  ► The dynamics along the $\beta$ direction forces the iterates to converge to the ground truth.

► **Explicit convergence rate** in 1-D case, $z = \min(|\beta|, |\beta^*|)$:
$$|M(\beta^*, \beta) - \text{sign}(\beta\beta^*)\beta^*| \leq \kappa(\beta^*, \beta, \sigma) \cdot |\beta - \text{sign}(\beta\beta^*)\beta^*| \qquad (1)$$

  ► Gaussian: $\kappa(\beta^*, \beta, \sigma) \leq \exp\left( -z^2/2\sigma^2 \right)$,

  ► Laplace: $\kappa(\beta^*, \beta, \sigma) \leq \frac{2\exp(-\frac{\sqrt{2}}{\sigma}z)}{1 + \exp(-2\frac{\sqrt{2}}{\sigma}z)}$

  ► Logistic: $\kappa(\beta^*, \beta, \sigma) \leq \frac{4\exp(-\frac{\pi z}{\sigma\sqrt{3}})}{1 + \exp(-\frac{2\pi z}{\sigma\sqrt{3}}) + 2\exp(-\frac{\pi z}{\sigma\sqrt{3}})}$
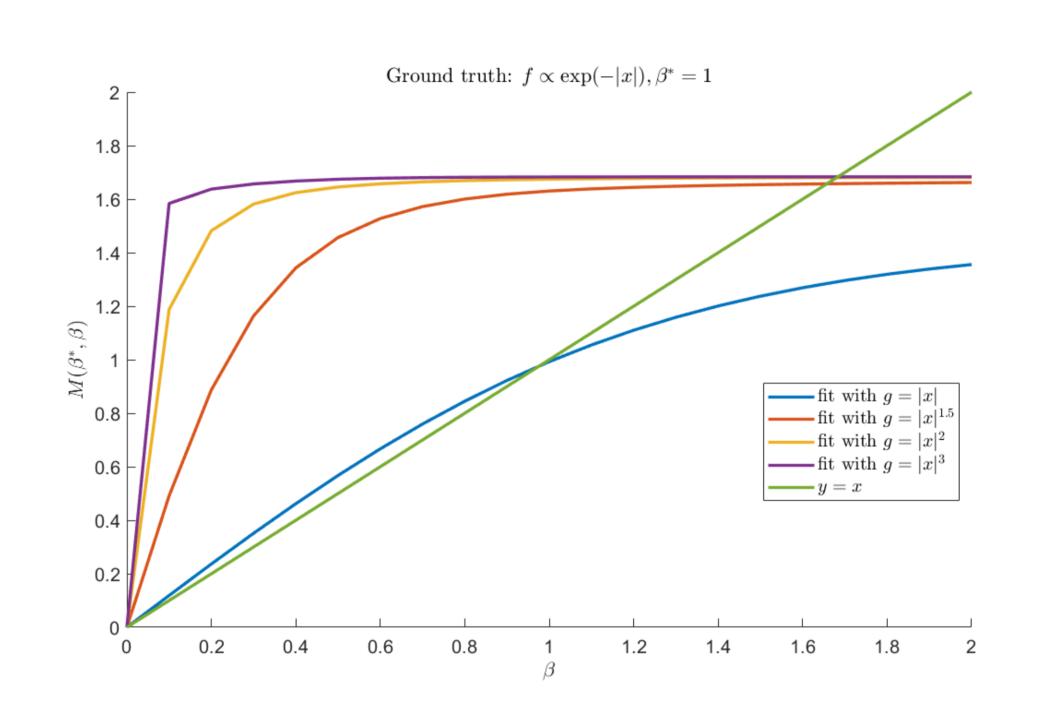
## Different Behaviors Compared to 2GMM

► In 1-D, the contraction in distance to the ground truth (1) holds for all $f \in \mathcal{F}$;

► In higher dimension, the contraction in $\ell_2$ distance still holds for Gaussian. However, there exists some log-concave distribution such that the $\ell_2$ distance strictly increases.

## Robustness under Model Mis-specification

In practice, we do not know $f$ that generates the data. Instead, we fit with some $\widehat{f} \in \mathcal{F}$. The LS-EM iterate under the mis-specification setting is:
$$\widehat{M(\beta^*, \beta)} = \mathbb{E}_{\boldsymbol{X} \sim D(\beta^*,\sigma)} \boldsymbol{X} \tanh\left( \frac{1}{2}\widehat{g}\left( \frac{1}{\sigma} \|\boldsymbol{X} + \beta\|_2 \right) - \frac{1}{2}\widehat{g}\left( \frac{1}{\sigma} \|\boldsymbol{X} - \beta\|_2 \right) \right).$$

► **Preserved properties**: two dimensional structure; angle decreasing property.

► **3-fixed points**: there are only 3 fixed points $\{\pm\overline{\beta}, 0\}$ in the direction of $\beta^*$.

► **Using Gaussian is a good choice:** when $\widehat{f}$ is Gaussian, $|\overline{\beta} - \beta^*| \leq 10\sigma$ if the SNR $|\beta^*|/\sigma$ is moderate.
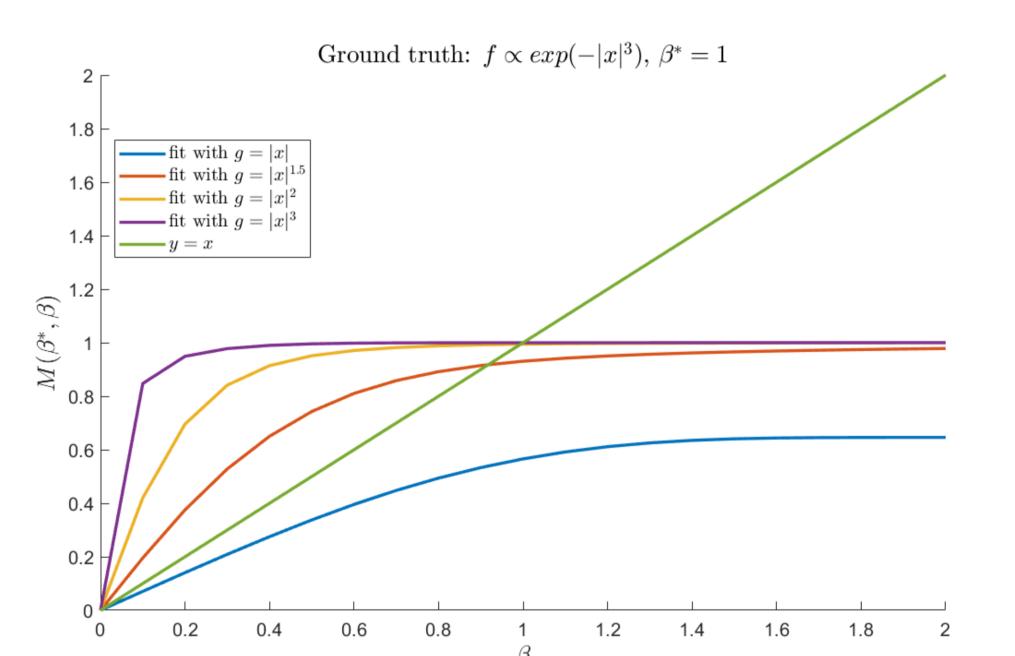


Figure: Ground truth: $f \propto \exp(-|x|)$, plot of $\widehat{M}(\beta^*, \beta)$ with $g = |x|, |x|^{1.5}, |x|^2$ and $|x|^3$.

Figure: Ground truth: $f \propto \exp(-|x|^3)$, plot of $\widehat{M}(\beta^*, \beta)$ with $g = |x|, |x|^{1.5}, |x|^2$ and $|x|^3$.

## Summary

► Two dimensional structure of the Least Squares EM under both correctly specified and mis-specified settings. Rotation invariance assumption guarantees this property;

► Angle decreasing of the Least Squares EM: convergence to the correct subspace. Log-concavity assumption and monotonicity of $g$ guarantee this property.